

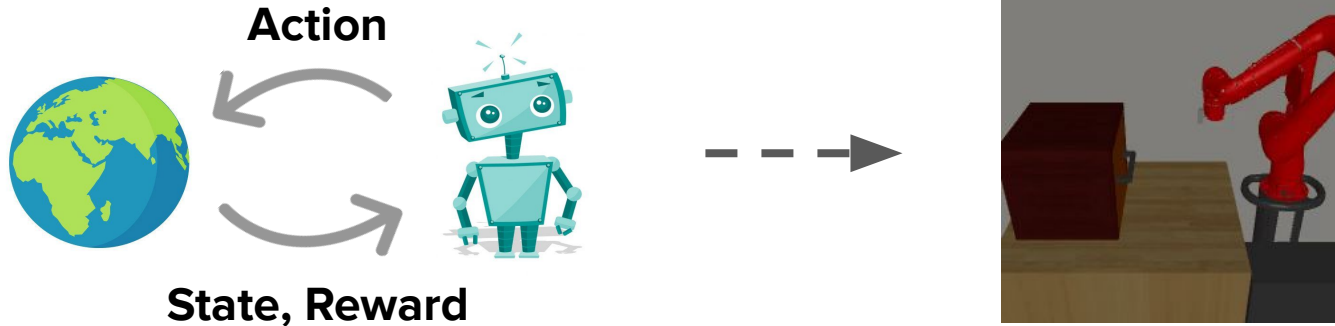
MURAL: Meta-Learning Uncertainty-Aware Rewards for Outcome-Driven Reinforcement Learning

Kevin Li*, Abhishek Gupta*, Ashwin Reddy, Vitchyr Pong, Aurick Zhou,
Justin Yu, Sergey Levine

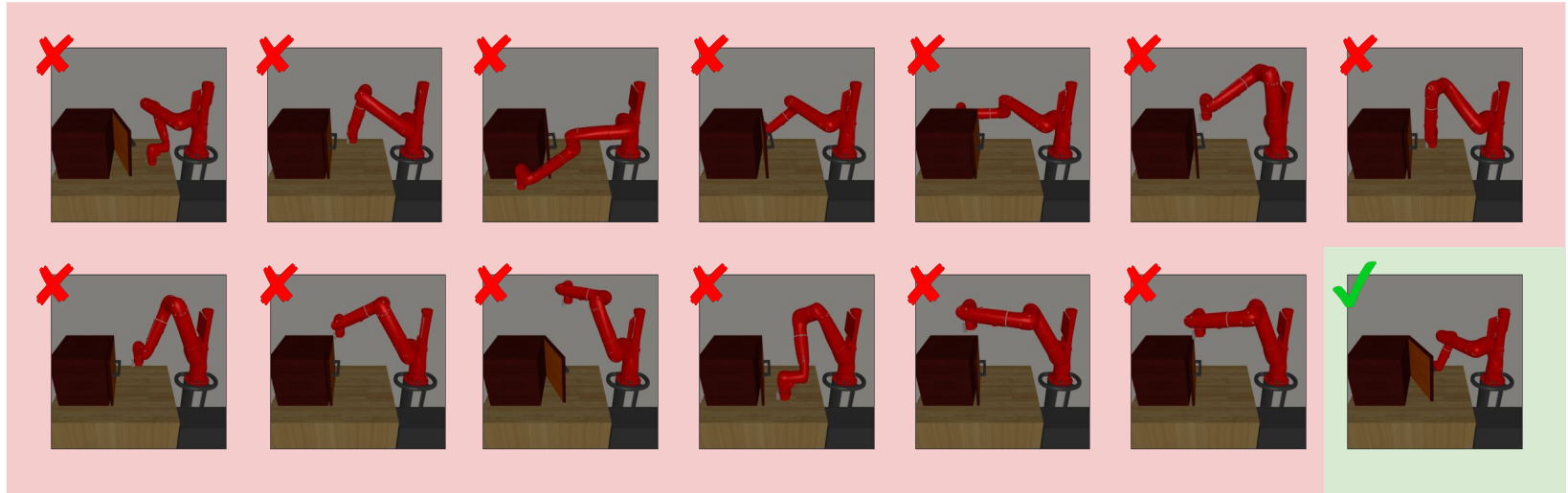


* = equal contribution

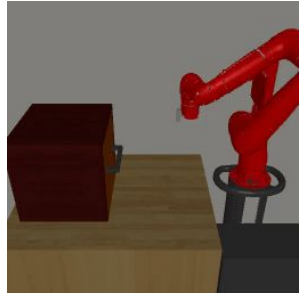
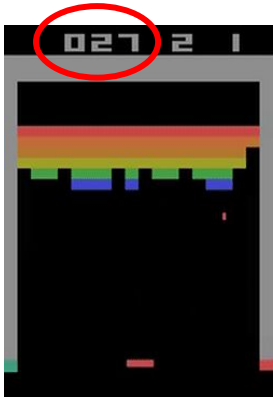
Reinforcement learning is a powerful framework for learning skills through interaction with the environment.



Often intractable — too many states to check!

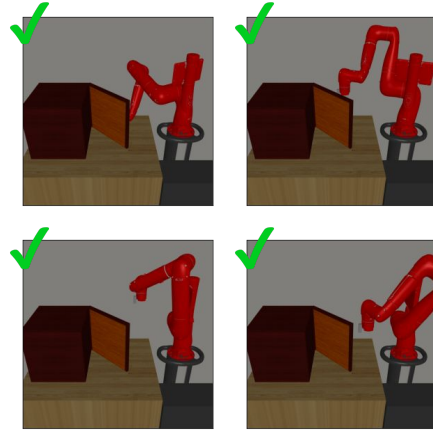


Reward shaping is helpful,
but not always possible...



???

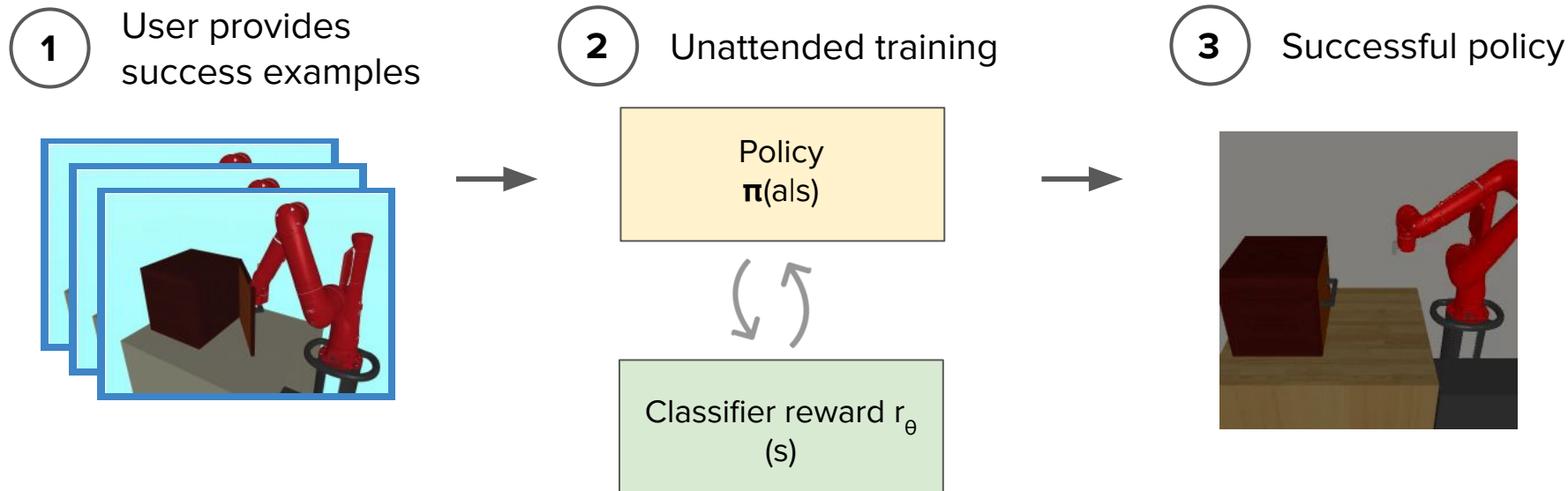
What if we formulate the RL problem
in terms of **successful outcomes**?



*Success: door is open at 45
degrees*

Classifier-Based Rewards

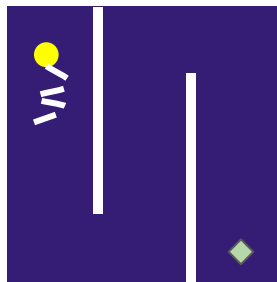
Idea: provide a set of desired outcomes, then train a classifier to distinguish between visited states and success states



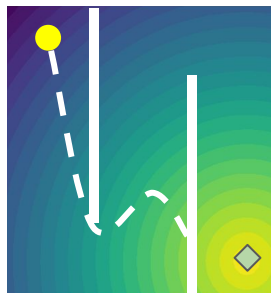
(Prior work: [Fu et al., 2018](#); [Singh et al., 2019](#), [Zhu et al., 2020](#))

Toy Example: 2D Maze

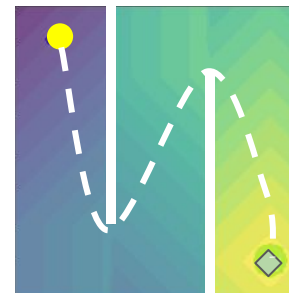
Sparse Reward
No reward signal



L2 Distance Reward
Misleading local optima

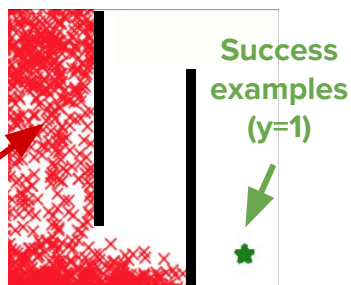


Ideal Shaped Reward
Guides agent toward goal

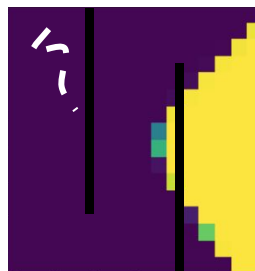


Possible
rewards

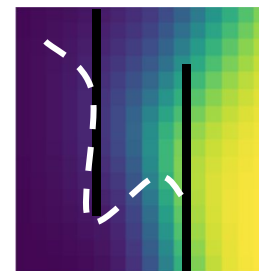
Training Set



Trained to convergence
Little reward signal



Early stopping + L2 reg
Misleading local optima



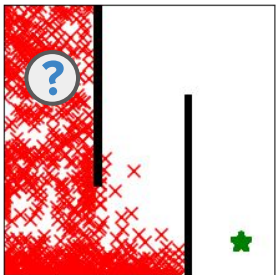
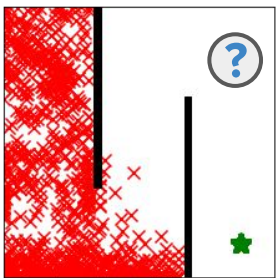
Standard
classifiers

Visited states
(y=0)

Success
examples
(y=1)

Key idea: use **Normalized Maximum Likelihood (NML)** to obtain reward classifiers that are uncertainty-aware and provide better shaping.

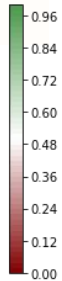
Query Point



y = 0



y = 1



Reward

$$\frac{1}{1 + 1} = 0.5$$

(high uncertainty)

$$\frac{0}{1 + 0} = 0$$

(same as original classifier)

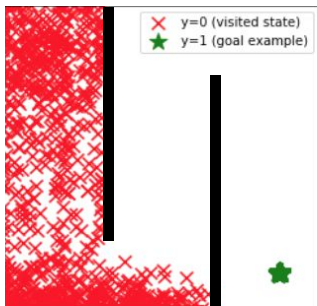
Key idea: use **Normalized Maximum Likelihood (NML)** to obtain reward classifiers that are uncertainty-aware and provide better shaping.

$$\theta_i = \arg \max_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mathcal{D} \cup (x_q, y=i)} [\log p_{\theta}(y|x)]$$

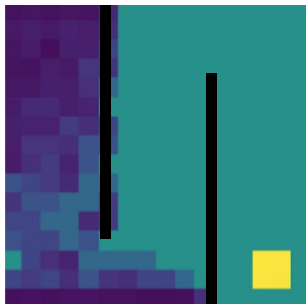
SLOW!!

$$p_{\text{CNML}}(y = i|x_q) = \frac{p_{\theta_i}(y = i|x_q)}{\sum_{j=1}^k p_{\theta_j}(y = j|x_q)}$$

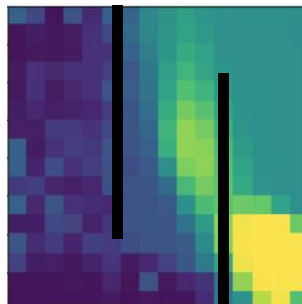
Training Set



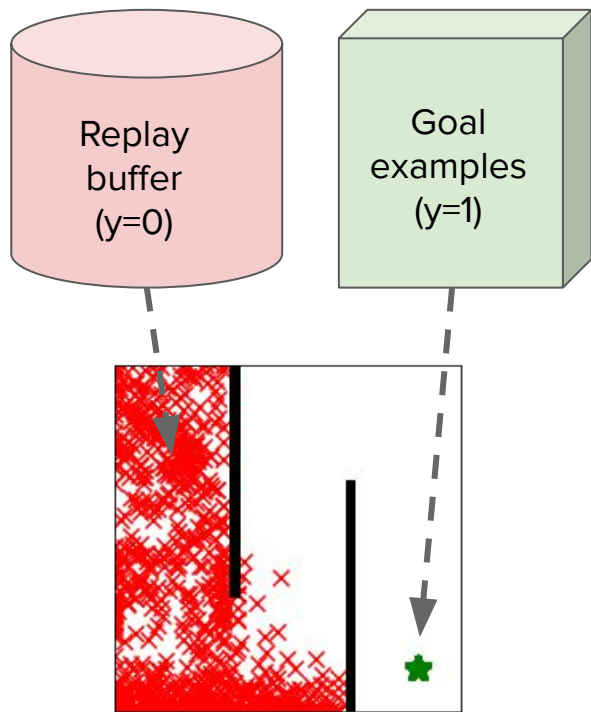
NML (discrete)



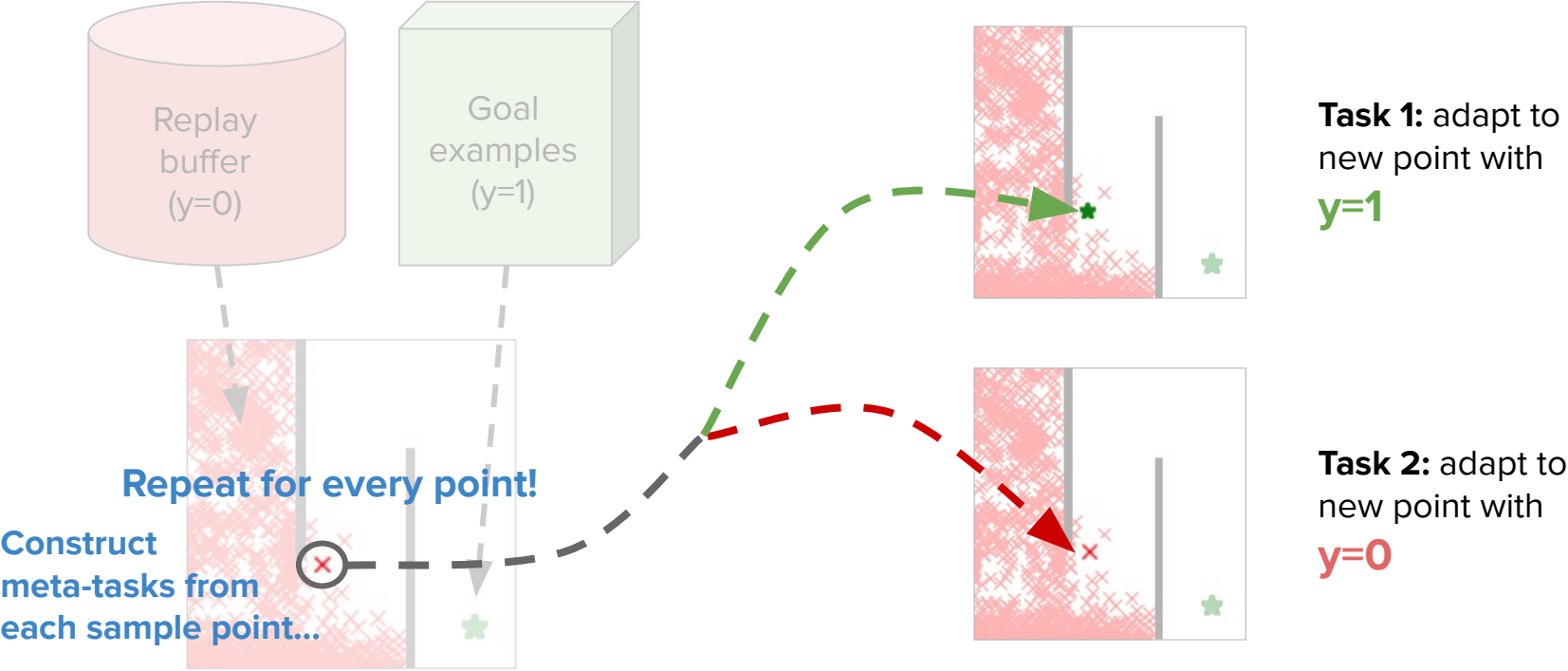
NML (with deep NNs)



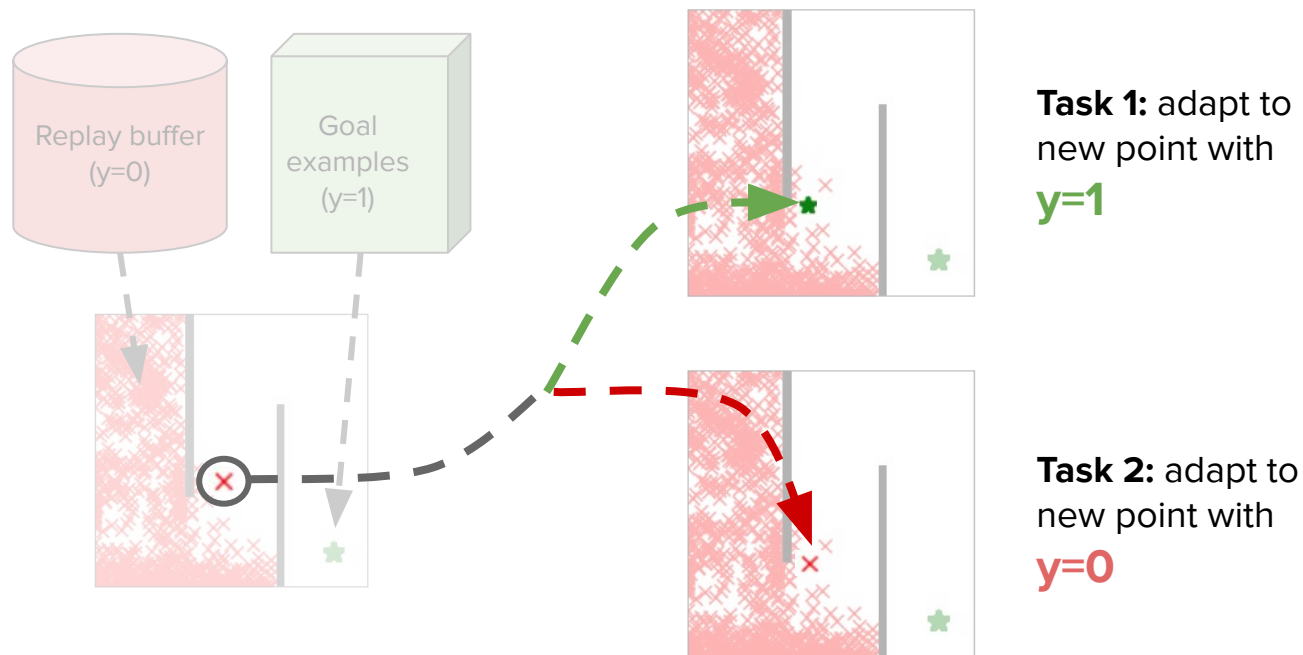
Meta-NML: efficient NML via meta-learning



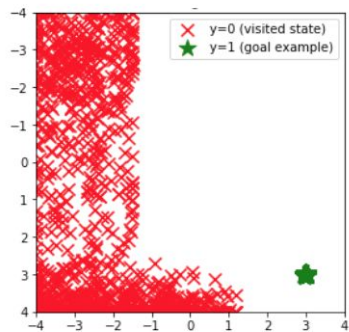
Meta-NML: efficient NML via meta-learning



Meta-NML: efficient NML via meta-learning



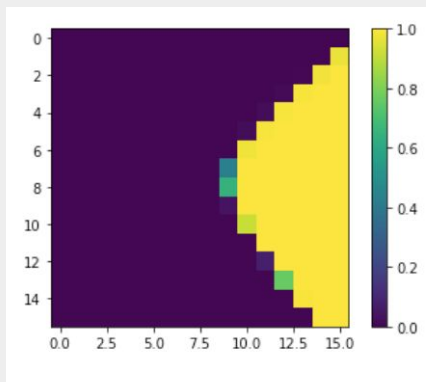
Training set



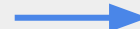
Meta-training



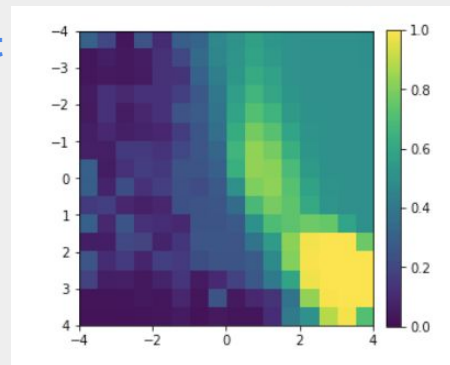
Model Initialization



1 gradient
step...



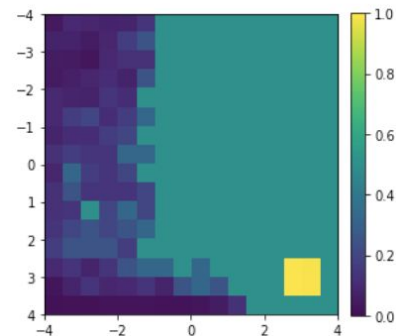
Meta-NML Rewards



Runtimes for a single epoch of RL

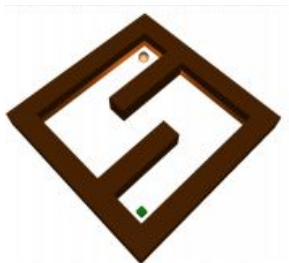
MLE Classifier	Meta-NML	Naive CNML
23.50s	39.05s	4hr 13min 34s

(Ideal NML Rewards)



Evaluation Domains

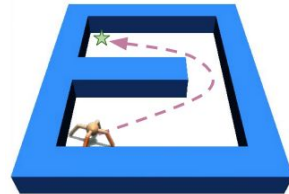
Navigation Tasks



Zigzag Maze

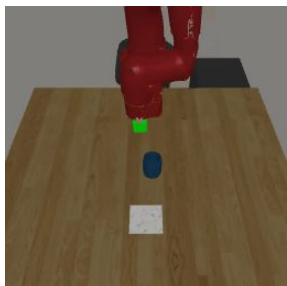


Spiral Maze

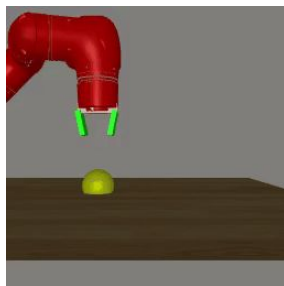


Ant Locomotion

Robotic Manipulation Tasks



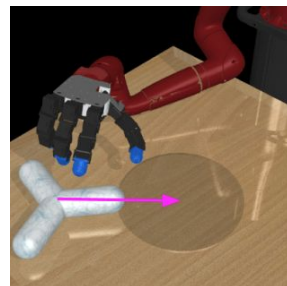
2D Pusher



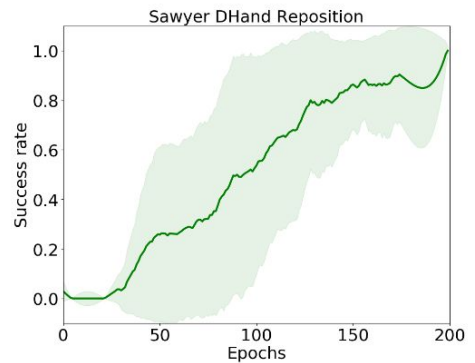
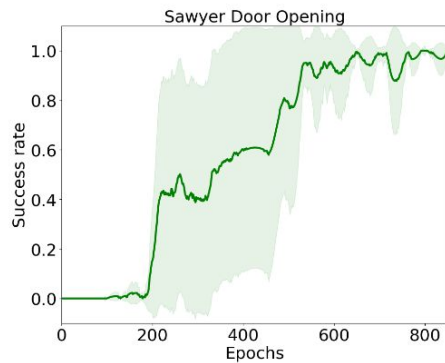
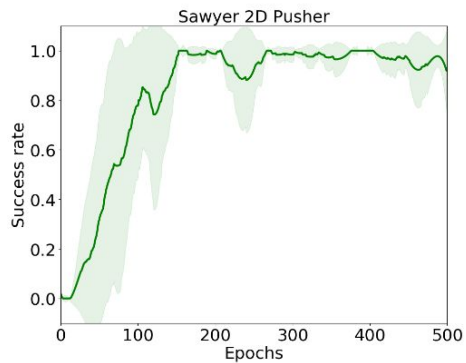
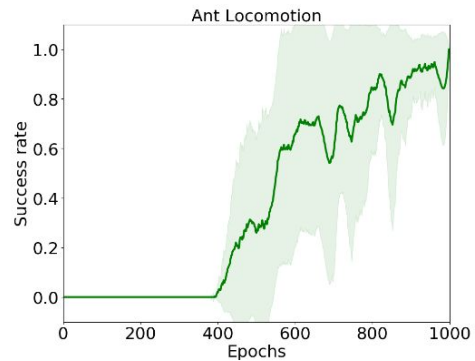
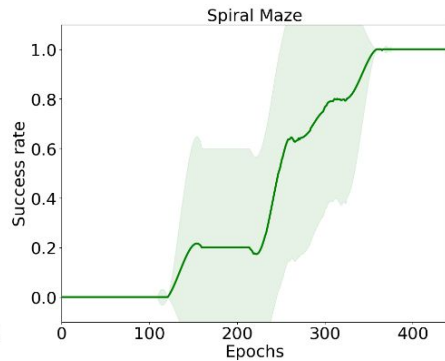
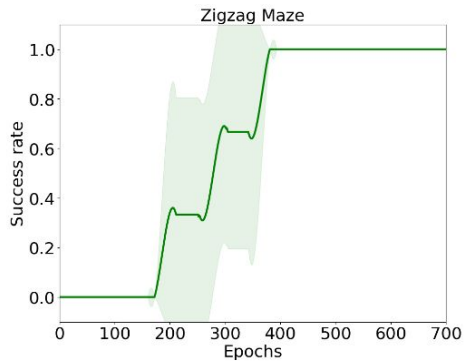
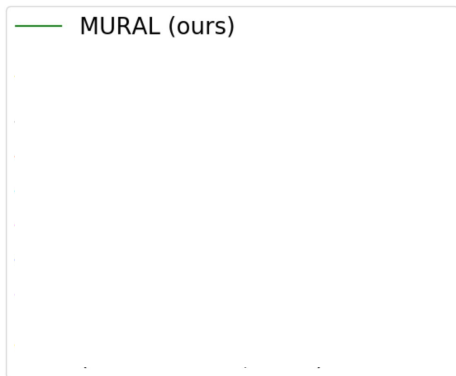
3D Pick-and-Place

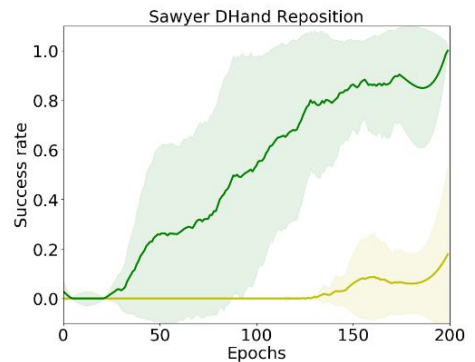
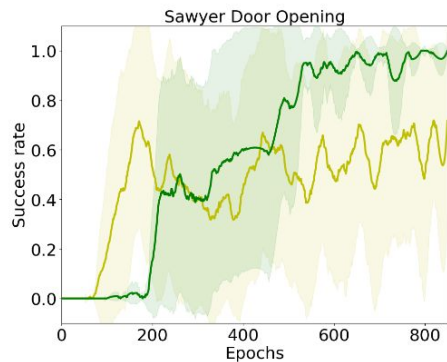
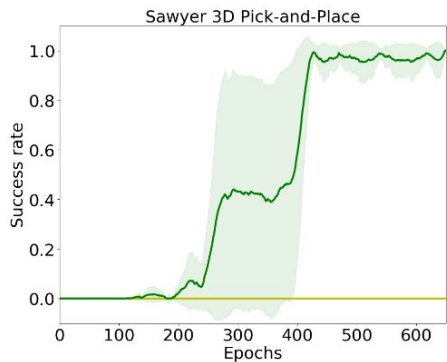
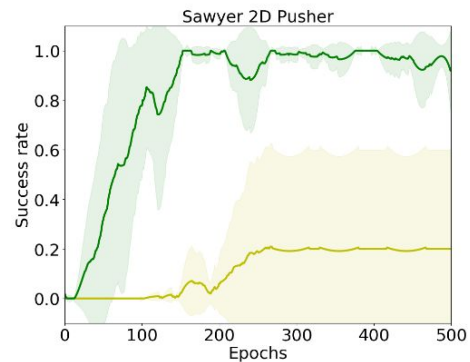
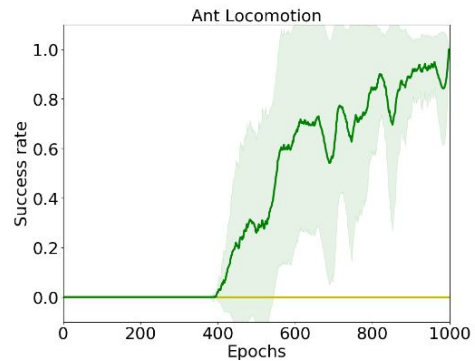
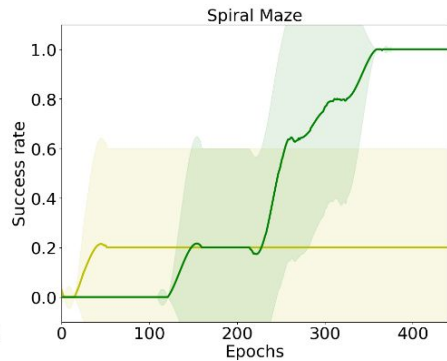
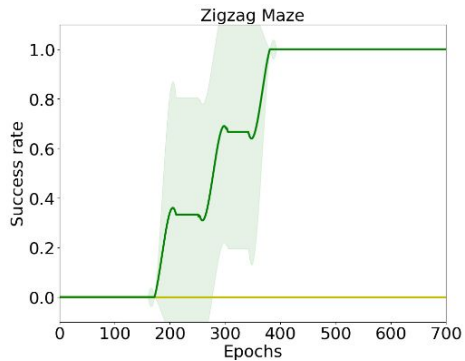


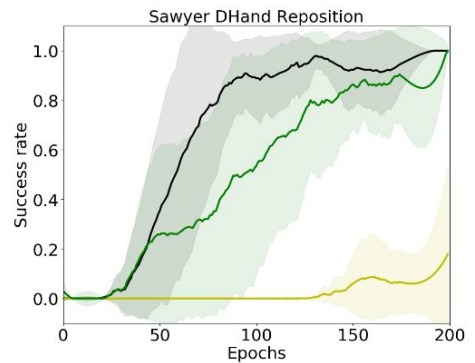
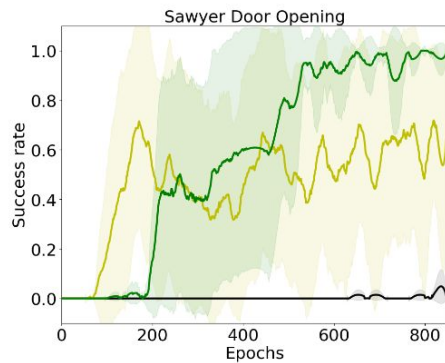
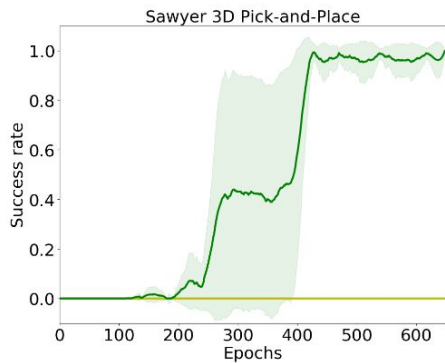
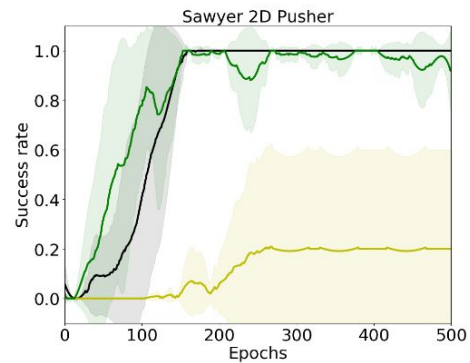
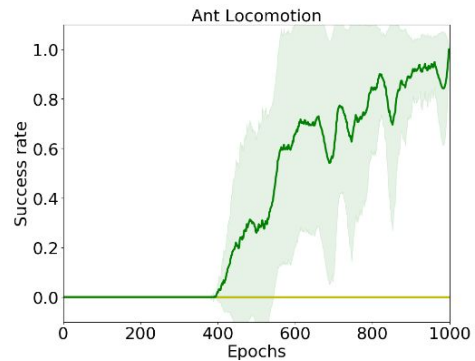
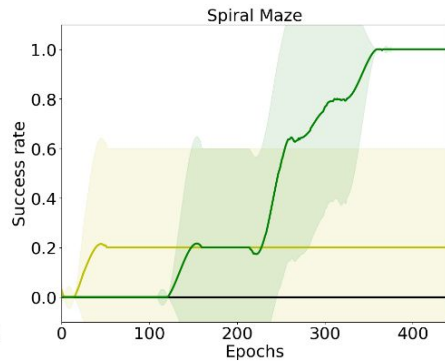
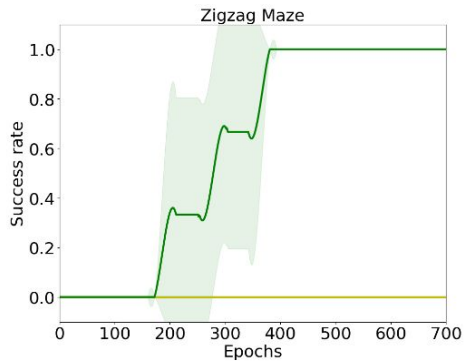
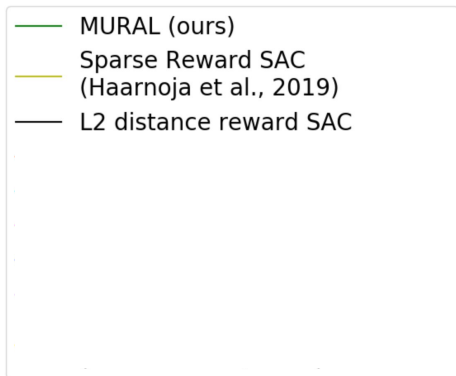
Door Opening

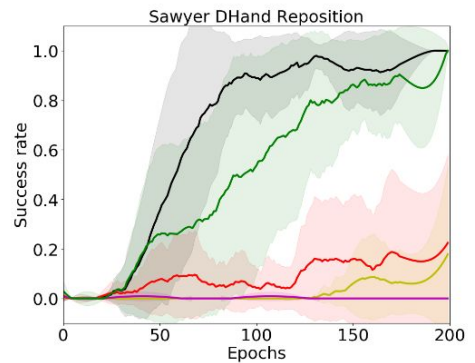
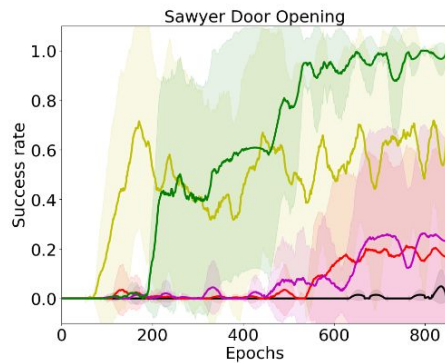
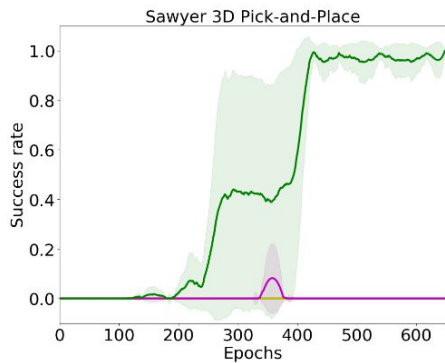
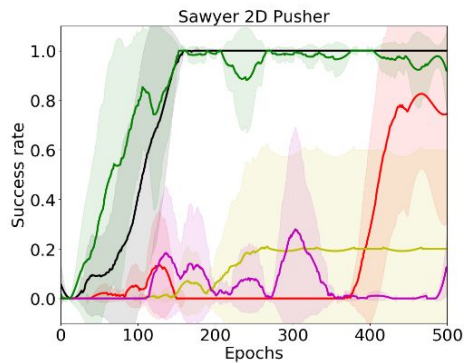
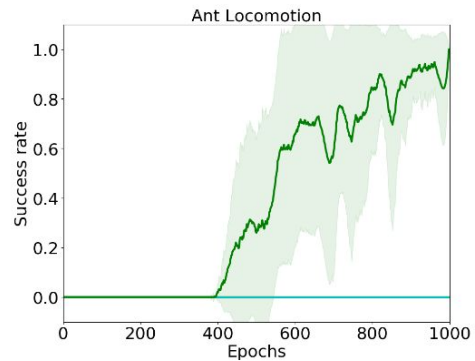
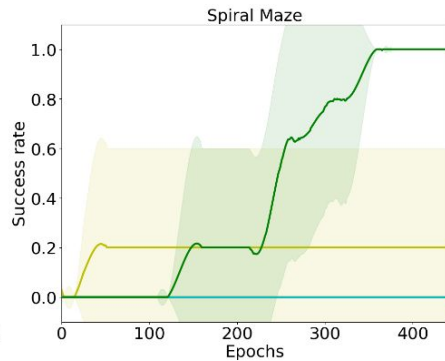
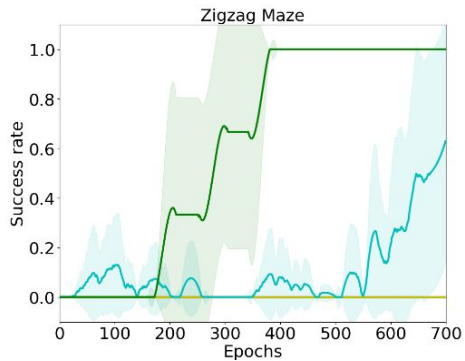
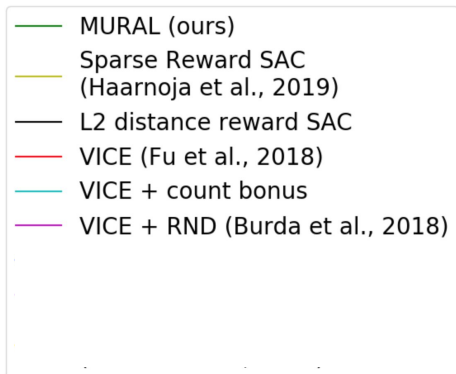


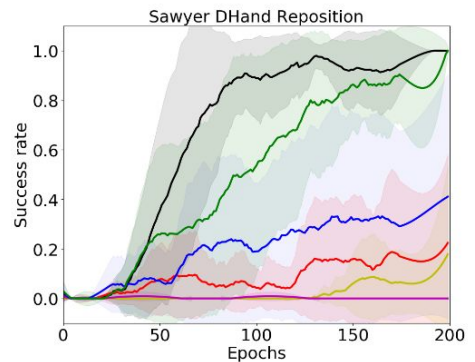
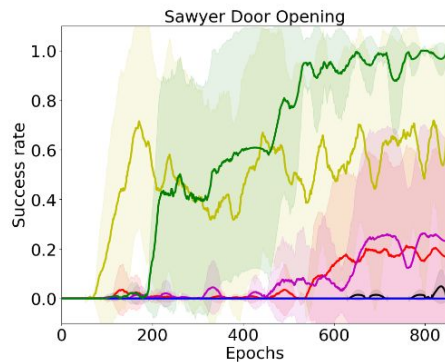
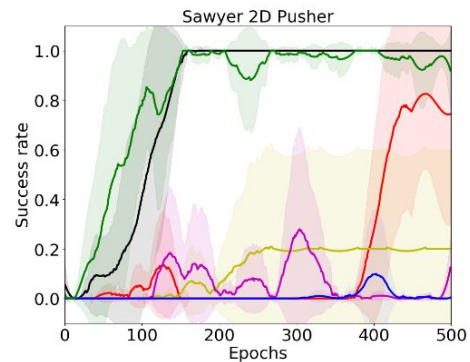
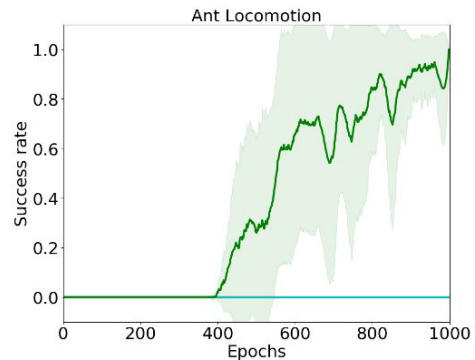
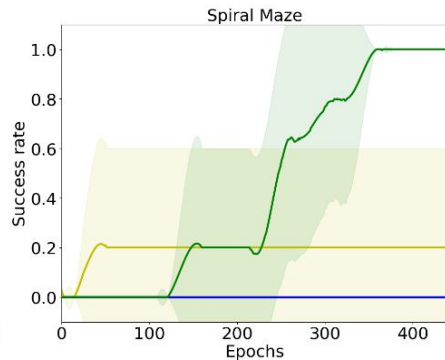
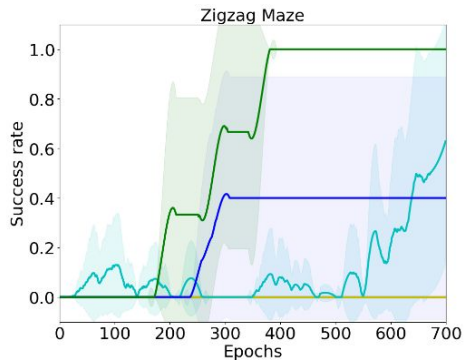
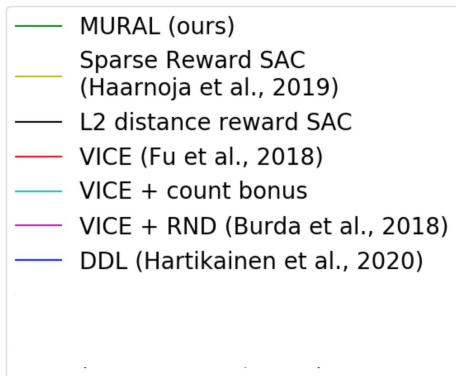
Dexterous Hand

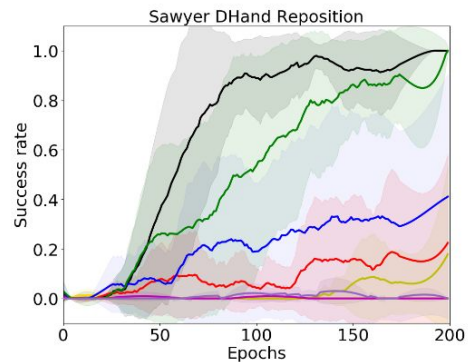
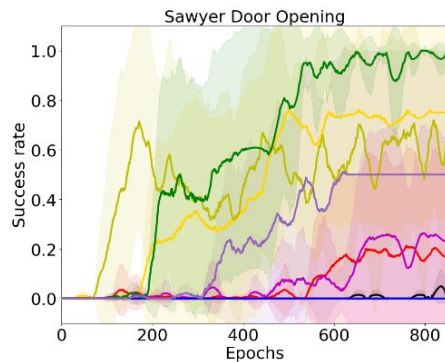
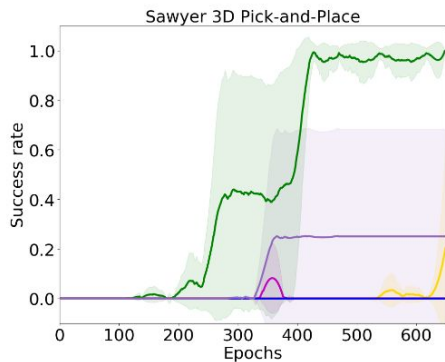
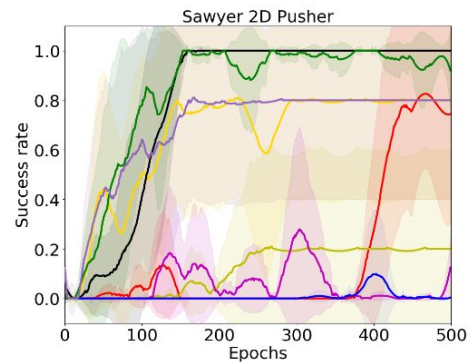
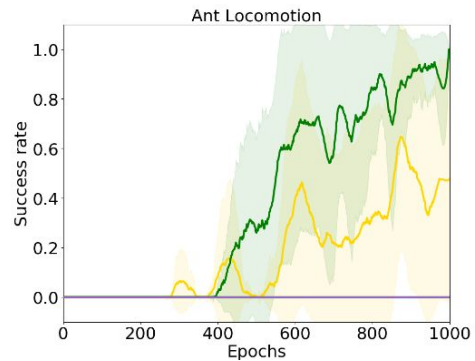
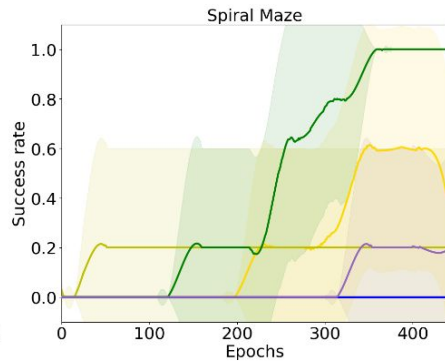
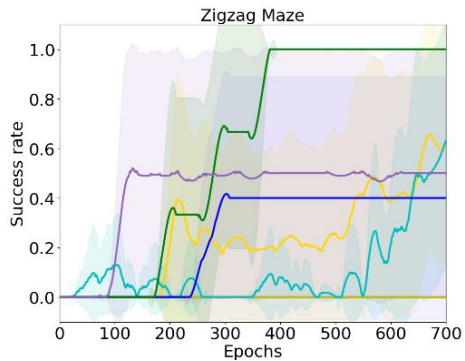
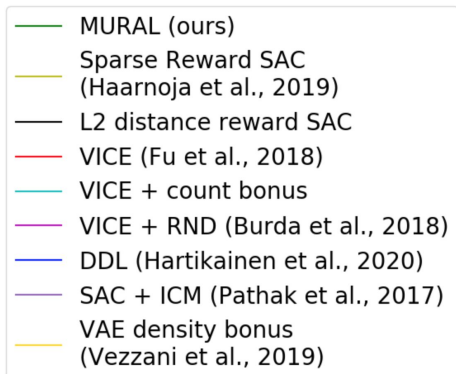












MURAL: Meta-Learning Uncertainty-Aware Rewards for Outcome-Driven Reinforcement Learning

Kevin Li*, Abhishek Gupta*, Ashwin Reddy, Vitchyr Pong, Aurick Zhou, Justin Yu, Sergey Levine



Website: <https://sites.google.com/view/mural-rl>

